

Recall $\vec{b} = (X^T X)^{-1} X^T \vec{y}$ was the result from $\arg\min_{\vec{w} \in \mathbb{R}^p} \|\vec{y} - X\vec{w}\|^2$ L1

This is not one of the estimation-generating procedures we talked about in 3A1.

Consider $\vec{\epsilon} \sim N_n(\vec{0}_n, \sigma^2 I_n)$, what is the MLE for $\vec{\beta}$?
 Our argument $\Rightarrow \vec{y} \sim N_n(X\vec{\beta}, \sigma^2 I_n) = f(\vec{y}; \vec{\beta}, \sigma^2, X) \Rightarrow$

$$\mathcal{L}(\vec{\beta}, \sigma^2; \vec{y}, X) = N_n(\vec{y} | X\vec{\beta}, \sigma^2 I_n) = \frac{1}{(2\pi)^{n/2} \sqrt{\det[\sigma^2 I_n]}} e^{-\frac{1}{2} (\vec{y} - X\vec{\beta})^T (\sigma^2 I_n)^{-1} (\vec{y} - X\vec{\beta})}$$

$$\mathcal{L}(\vec{\beta}, \sigma^2; \vec{y}, X) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \|\vec{y} - X\vec{\beta}\|^2$$

$$\frac{\partial}{\partial \vec{\beta}} [\mathcal{L}] = -\frac{1}{2\sigma^2} \frac{\partial}{\partial \vec{\beta}} [\|\vec{y} - X\vec{\beta}\|^2] \stackrel{\text{set}}{=} 0 \Rightarrow \frac{\partial}{\partial \vec{\beta}} \left[\sum (\vec{y}_i - \vec{x}_i \vec{\beta})^2 \right] = 0$$

This is exactly the same problem solved before! $\Rightarrow \hat{\vec{\beta}}_{MLE} = (X^T X)^{-1} X^T \vec{y} = \vec{b}$,
 i.e. the least squares estimate is the MLE for $\vec{\beta}$!

Note: this is the MLE regardless of the value of σ^2 !

Analogue: \bar{X} has $\hat{\theta}_{MLE}$ for $X_i \sim \mathcal{N}(\theta, \sigma^2)$ regardless of value of σ^2 .

$$\frac{\partial}{\partial \sigma^2} [\mathcal{L}] = -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \|\vec{y} - X\vec{\beta}\|^2 \stackrel{\text{set}}{=} 0 \Rightarrow -n + \frac{1}{\sigma^2} \|\vec{y} - X\vec{\beta}\|^2 = 0$$

$$\Rightarrow \hat{\sigma^2}^{MLE} = \frac{1}{n} \|\vec{y} - X\vec{\beta}\|^2 \quad \checkmark \quad \text{Substitute } \vec{\beta} = \hat{\vec{\beta}}_{MLE}$$

$$= \frac{1}{n} \|(\mathbf{I} - H) \vec{y}\|^2$$

$$= \frac{1}{n} \|\vec{e}\|^2 = \frac{1}{n} \sum e_i^2 = \frac{1}{n} SSE$$

Oh Hw!!

Remember in 3A2W, if $p+1 \geq n$, OLS fails since $X^T X$ is no longer invertible. But what if we still want inference?

Bayesian shrinkage estimators can help us.

Let's assume a prior on $\vec{\beta}$: $\beta_0, \beta_1, \dots, \beta_p \stackrel{iid}{\sim} N(0, \sigma^2)$

$\Rightarrow f(\vec{\beta}) = N_{p+1}(\vec{0}_{p+1}, \tau^2 I_{p+1})$. If $\tau^2 \rightarrow \infty$ this prior is the Laplace shrink prior for all β_j

And also on σ^2
 $f(\sigma^2) \propto \frac{1}{\sigma^2} \propto \text{Inverse Gamma}(0, 0)$ and independent of $\vec{\beta}$

$$\begin{aligned} f(\vec{\beta}, \sigma^2 | X, \vec{y}) &\propto f(\vec{y} | \vec{\beta}, \sigma^2, X) f(\vec{\beta}, \sigma^2 | X) \\ &= \frac{1}{(2\pi)^{n/2} \sqrt{\det[\sigma^2 I_n]}} e^{-\frac{1}{2} (\vec{y} - X\vec{\beta})^T (\sigma^2 I_n)^{-1} (\vec{y} - X\vec{\beta})} \\ &\quad \frac{1}{(2\pi)^{\frac{p+1}{2}} \sqrt{\det[\tau^2 I_{p+1}]}} e^{-\frac{1}{2} (\vec{\beta} - \vec{0}_{p+1})^T (\tau^2 I_{p+1})^{-1} (\vec{\beta} - \vec{0}_{p+1})} \cdot \frac{1}{\sigma^2} \\ &\propto (\sigma^2)^{-\frac{n}{2}-1} e^{-\frac{1}{2\sigma^2} \|\vec{y} - X\vec{\beta}\|^2} e^{-\frac{1}{2\tau^2} \|\vec{\beta}\|^2} \end{aligned}$$

Since we care currently about inference for $\vec{\beta}$, let's find the marginal posterior:

$$\begin{aligned} f(\vec{\beta} | X, \vec{y}) &= \int_0^\infty f(\vec{\beta}, \sigma^2 | X, \vec{y}) d\sigma^2 \propto e^{-\frac{1}{2\tau^2} \|\vec{\beta}\|^2} \int_0^\infty (\sigma^2)^{-\frac{n}{2}-1} e^{-\frac{\|\vec{y} - X\vec{\beta}\|^2/2}{\sigma^2}} d\sigma^2 \\ &= e^{-\frac{1}{2\tau^2} \|\vec{\beta}\|^2} \frac{\sqrt{\frac{n}{2}}}{(\|\vec{y} - X\vec{\beta}\|^2/2)^{n/2}} \propto e^{-\frac{1}{2\tau^2} \|\vec{\beta}\|^2} \frac{1}{\|\vec{y} - X\vec{\beta}\|^n} \end{aligned}$$

Kernel for inversegamma(a, b)
where $a = \frac{n}{2}$, $b = \|\vec{y} - X\vec{\beta}\|^2/2$

This is not the result of a known kernel. To do inference, we could Gibbs sample it. But if we just care about estimation, we can calculate the $\hat{\beta}_{MAP}$

$$\hat{\beta}_{MAP} := \underset{\beta}{\operatorname{argmax}} \{ f(\beta, \sigma^2 | x, y) \} = \underset{\beta}{\operatorname{argmax}} \{ \ln(f(\beta, \sigma^2 | x, y)) \}$$

$$= \underset{\beta}{\operatorname{argmax}} \left\{ \underbrace{\left(-\frac{1}{2}\right) \ln(\sigma^2)}_{\text{constant}} - \frac{1}{2\sigma^2} \|\tilde{y} - X\beta\|^2 - \frac{1}{2\sigma^2} \|\beta\|^2 \right\}$$

$$= \underset{\beta}{\operatorname{argmax}} \left\{ \underbrace{-\frac{1}{2\sigma^2}}_{\text{constant}} \left(\|\tilde{y} - X\beta\|^2 + \frac{\sigma^2}{\tau^2} \|\beta\|^2 \right) \right\}$$

$$= \underset{\beta}{\operatorname{argmin}} \left\{ \underbrace{\|\tilde{y} - X\beta\|^2}_{\text{residual}} + \underbrace{\frac{\sigma^2}{\tau^2} \|\beta\|^2}_{\text{penalty}} \right\}$$

Now if we just care about estimation technique a/g 392w....

$$\downarrow \text{let } \lambda = \frac{\sigma^2}{\tau^2} \geq 0 \quad = 0 \text{ if } \tau^2 \rightarrow \infty$$

$$= \underset{\beta}{\operatorname{argmin}} \{ SSE + \lambda \|\beta\|^2 \} = \text{Ridge Regression Estimator}$$

hyperparameter you choose. If $\lambda = 0$

$$\Rightarrow \beta = OLS.$$

If λ is large $\beta \approx \vec{0}_{p+1}$

which forces the model to shrink β_j 's towards $\vec{0}$ β_j 's that are very useful for prediction.

$$\hat{y}_{ridge} = X \hat{\beta}_{ridge} \in \operatorname{Colsp}(X) \text{ but it is not the orthogonal projection}$$

This process is known as "regularization" as it's making the $\hat{\beta}$ estimates more "regular" i.e. simpler, closer to $\vec{0}$

see next page

$$\hat{\beta}_{ridge} = (X^T X + \lambda I_{p+1})^{-1} X^T \tilde{y}$$

$\hat{\beta}_{ridge}$

$$\begin{aligned}
 &= \frac{\partial}{\partial \vec{\beta}} \left[(\vec{y} - X\vec{\beta})^T (\vec{y} - X\vec{\beta}) + \lambda \vec{\beta}^T \vec{\beta} \right] \\
 &= \frac{\partial}{\partial \vec{\beta}} \left[\vec{y}^T \vec{y} - 2\vec{\beta}^T X^T \vec{y} + \vec{\beta}^T X^T X \vec{\beta} + \vec{\beta}^T (\lambda I_{p+1}) \vec{\beta} \right] \\
 &= \frac{\partial}{\partial \vec{\beta}} \left[-2\vec{\beta}^T X^T \vec{y} + \vec{\beta}^T (X^T X + \lambda I_{p+1}) \vec{\beta} \right] \\
 &= -2X^T \vec{y} + 2(X^T X + \lambda I_{p+1}) \vec{\beta} \stackrel{\text{set}}{=} 0 \text{ to find min.} \\
 \Rightarrow X^T \vec{y} &= (X^T X + \lambda I_{p+1}) \vec{\beta} \Rightarrow \vec{b}_{\text{ridge}} = (X^T X + \lambda I_{p+1})^{-1} X^T \vec{y}
 \end{aligned}$$

always invertible
 if $\lambda > 0$
 as it forces full rank.

Recall that OLS was an unbiased estimate:

$$\begin{aligned}
 \vec{B} &= (X^T X)^{-1} X^T \vec{y}, \quad E(\vec{y}) = E[X(\beta + \vec{\epsilon})] = X\beta + E(\vec{\epsilon}) = X\beta + \vec{0}_n = X\beta \\
 E(\vec{B}) &= E[\vec{y}] = (X^T X)^{-1} X^T E(\vec{y}) = (X^T X)^{-1} X^T X \beta = \beta
 \end{aligned}$$

Ridge is biased:

$$E(\vec{B}_{\text{ridge}}) = (X^T X + \lambda I_{p+1})^{-1} X^T X \beta \quad \text{as } \lambda \rightarrow 0 \Rightarrow \text{bias} \rightarrow 0.$$

But it compensates for this bias with lower variance!

Hence the total MSE is usually lower!